

LEVEL

AD47496

UCLA-ENG-7983
MAY 1979

ARPA CONTRACT NO. DAHC-15-C-0368

**ADVANCED TELEPROCESSING SYSTEMS
FINAL TECHNICAL REPORT**

May 31, 1979

DTIC

MAR 11 1980

Principal Investigator: Leonard Kleinrock
Co-Principal Investigators: Gerald Estrin
Richard R. Muntz
Gerald Popek

**COMPUTER SYSTEMS
MODELING AND ANALYSIS GROUP**

ADA081916

This document has been approved
for public release and sale; its
distribution is unlimited.

COMPUTER SCIENCE DEPARTMENT

School of Engineering and Applied Science
University of California
Los Angeles



DDC FILE COPY

80 3 05 037

UNCLASSIFIED

SECURITY CLASSIFICATION OF THIS PAGE (When Data Entered)

REPORT DOCUMENTATION PAGE		READ INSTRUCTIONS BEFORE COMPLETING FORM
1. REPORT NUMBER 14 UCLA-ENG-7983	2. GOVT ACCESSION NO.	3. RECIPIENT'S CATALOG NUMBER
4. TITLE (and Subtitle) 6 ADVANCED TELEPROCESSING SYSTEMS FINAL TECHNICAL REPORT		5. TYPE OF REPORT & PERIOD COVERED Final Technical Report, 12/1/75 - 3/31/79
7. AUTHOR(s) 10 Leonard/Kleinrock Richard R. Muntz Gerald/Estrin Gerald/PopeK 15		8. CONTRACT OR GRANT NUMBER(s) DAHC 15-73-C-0368 WARPA Order-2496
9. PERFORMING ORGANIZATION NAME AND ADDRESS School of Engineering and Applied Science University of California Los Angeles, CA 90024		10. PROGRAM ELEMENT, PROJECT, TASK AREA & WORK UNIT NUMBERS ARPA Order No. 2496
11. CONTROLLING OFFICE NAME AND ADDRESS Advanced Research Projects Agency 1400 Wilson Boulevard Arlington, VA 22209		12. REPORT DATE May 31, 1979
14. MONITORING AGENCY NAME & ADDRESS (if different from Controlling Office) 9 Final technical rept. 1 Dec 75 - 31 Mar 79		13. NUMBER OF PAGES 17
16. DISTRIBUTION STATEMENT (of this Report) Approved for public release; distribution unlimited 11 31 May 79		15. SECURITY CLASS. (of this report) Unclassified
17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different from Report)		18a. DECLASSIFICATION/DOWNGRADING SCHEDULE
18. SUPPLEMENTARY NOTES		
19. KEY WORDS (Continue on reverse side if necessary and identify by block number)		
20. ABSTRACT (Continue on reverse side if necessary and identify by block number) This is the final report for the ARPA Advanced Teleprocessing Systems Contract DAHC-15-C-0368 by the UCLA Computer Science Department covering the period from December 1, 1975 to March 31, 1979. During this contract period we were engaged in three principal tasks: 405749 THIS REPORT COVERS THE FOLLOWING TASKS (over) 7 50B		

UNCLASSIFIED

SECURITY CLASSIFICATION OF THIS PAGE(When Data Entered)

(20)

- (1) Packet Radio Studies. Analytical evaluation of random multi-access modes for single and multiple repeater environments, networking issues (i.e., multiple channels, repeater range, configuration, etc.), dynamic routing procedures and network protocols; experimental evaluation of the packet radio network system (measurement and simulation),
- (2) Satellite Communications. Experimental evaluation of the satellite network using SIMPS in the Atlantic experiment. Design of experiments for evaluation and comparison of channel access and control strategies. Specification and implementation of data reduction algorithms and experiments,
- (3) Resource allocation and sharing. Development of scheduling and resource allocation strategies for packet communications.

Included is a copy of the paper by the Principal Investigator (Leonard Kleinrock), "Principles and Lessons in Packet Communications," which appeared in the the Proceedings of the IEEE, November 1978, and which is an excellent summary of the activities and results of the research conducted during this contract period. Also included is a complete bibliography of all published works which were supported during this research contract.

Accession For	
NTIS GR&I	
DGC TAB	
Unannounced	
Justification	
By	
Distribution/	
Availability Codes	
Dist.	Avail and/or special
A	

UNCLASSIFIED

SECURITY CLASSIFICATION OF THIS PAGE(When Data Entered)

UCLA-ENG-7983
May 1979

ADVANCED RESEARCH PROJECTS AGENCY

ADVANCED TELEPROCESSING SYSTEMS:

ACCESS, ALLOCATION AND CONTROL

Final Technical Report

May 31, 1979

ARPA Contract: DAHC 15-73-C-0368
ARPA Order No.: 2596

Principal Investigator: Leonard Kleinrock
(213) 825-2543

Co-Principal Investigators: Gerald Estrin
Richard R. Muntz
Gerald Popek

Contact Period: June 15, 1973 - March 31, 1979
Report Period: December 1, 1975 - March 31, 1979

Computer Science Department
School of Engineering and Applied Science
University of California, Los Angeles

ADVANCED TELEPROCESSING SYSTEMS: ACCESS, ALLOCATION AND CONTROL

Advanced Research Projects Agency Final Technical Report

May 31, 1979

This final report covers the period from December 1, 1975 to March 31, 1979 for ARPA Contract Number DAHC-15-73-C-0368. The research conducted during this period has been amply reported upon in our Semi-Annual Technical Reports as well as in the published professional literature. Consequently the format of this report is quite simple and brief as have been our previous final technical reports. Below we restate the research tasks which have been the subject of this research project. We follow that with an extensive bibliography of those published papers which have appeared during the period from December 1, 1975 to June 30, 1977 plus one specific paper and two dissertations that appeared between July 1, 1977 and March 31, 1979. Although this contract formally ended on March 31, 1979, the period from July 1, 1977 beyond was simply a no-cost extension of funds to complete and publish studies that were begun during the main contract period. (See our final technical report, UCLA-ENG-7696, dated November 30, 1975 for a list of publications during the period June 15, 1973 to November 30, 1975.) At the end of this final report we reproduce a copy of the paper "Principles and Lessons in Packet Communications", which was authored by the Principal Investigator (L. Kleinrock) and which was published in the Proceedings of the IEEE, November 1978, pp. 1320-1329. This latter paper is a summary of the research activities and the research results of this contract, thereby providing an excellent summary and final report.

During this contract period we have been engaged in the following three principal tasks (a fourth task on "Secure Systems" is currently under the direction of Professor Gerald Popek and is reported upon by him in other documents):

Task I. - Packet Radio Studies

- A. Analytical evaluation of random multi-access modes for single and multiple repeater environments, networking issues, (i.e., multiple channels repeater range, configuration, etc.), dynamic routing procedures and network protocols.
- B. Experimental evaluation of the packet radio network system (i.e., measurement and simulation).

Task II. - Satellite Communications

Experimental evaluation of the satellite network using SIMPS in the Atlantic experiment. Design of experiments for evaluation and comparison of channel access and control strategies. Specification and implementation of data reduction algorithms and experiments.

Task III. - Resource Allocation and Sharing

Development of scheduling and resource allocation strategies for packet communications.

These tasks have all advanced significantly during this period. The output of this research in the period referred to in this final report takes the form of five Ph.D dissertations, nineteen publications of the principal faculty, staff and students which have appeared in the professional literature, plus fifteen Packet Radio Temporary Notes and Packet Satellite Program Working Notes (this does not include the twelve publications listed below which have come out of the task in "Secure Systems"). These publications are listed in the bibliography below. That bibliography is grouped into research areas as follows:

1. Resource Allocation and Sharing
2. Security Systems
3. Satellite Communications
4. Packet Radio Studies
5. Dissertations

Thus we offer as one measure of our achievements during this contract period the publications listed below. In addition to these publications there have been a large number of presentations at professional conferences around the world. The substance of the lessons and principles that have been established will be found in the paper following the publications list.

Bibliography

Resource Allocation and Sharing

Kamoun, F. and L. Kleinrock, "Analysis of Shared Storage in a Computer Network Environment," *Proceedings of the Ninth Hawaii International Conference on System Sciences*, Honolulu, January 1976, pp. 89-92.

Kleinrock, L., W. Naylor and H. Opderbeck, "A Study of Line Overhead in the ARPANET," *Communications of the ACM*, Vol. 19, January 1976, pp. 3-12.

DeWitt, H. and M. Krieger, "Expected Structure of Euclidean Graphs," Symposium on New Directions and Recent Results in Algorithms and Complexity, Carnegie-Mellon University, April 1976.

Krieger, M., "Random Graphs and Minimal Spanning Tree Algorithms," Symposium on New Directions and Recent Results in Algorithms and Complexity, Carnegie Mellon University, April 1976.

Kleinrock, L., "ARPANET Lessons," *Conference Record, International Conference on Communications*, Philadelphia, Pennsylvania, June 1976, pp. 20-1 to 20-6.

Kleinrock, L. and F. Kamoun, "Data Communications Through Large Packet-Switching Networks," *Proceedings of the Eighth International Teletraffic Congress*, Sydney, Australia, November 1976, pp. 521-1 to 521-10.

Naylor, W. and L. Kleinrock, "On Effect of Periodic Routing Updates in Packet-Switched Networks," *Conference Record, National Telecommunications Conference*, Dallas, Texas, November 1976, pp. 16.2-1 to 16.2-7.

Kleinrock, L. "On Communications and Networks," *IEEE Transactions on Computers*, 25th Anniversary Issue, Vol. C-25, December 1976, pp. 1326-1335.

Gerla, M. and L. Kleinrock, "On the Topological Design of Distributed Computer Networks," *IEEE Transactions on Communications*, Vol. COM-25, January 1977, pp. 48-60.

Kleinrock, L. and H. Opderbeck, "Throughput in the ARPANET - Protocols and Measurement," *IEEE Transactions on Communications*, Vol. COM-25, January 1977, pp. 95-104.

Kleinrock, L. and F. Kamoun, "Hierarchical Routing for Large Networks, Performance Evaluation and Optimization," *Computer Networks*, Vol. 1, No. 3, January 1977, pp. 155-174.

Kermani, P. and L. Kleinrock, "Analysis of Buffer Allocation Schemes in a Multiplexing Node," *Conference Record, International Conference on Communications*, Vol. 2, Chicago, Illinois, June 1977, pp. 30.4-269 to 275.

Kleinrock, L., and P. Kermani, "A Network Algebra for the Performance Evaluation of Interconnected Networks," *Conference Record, National Telecommunications Conference*, Birmingham, Alabama, December 1978, pp. 44.1.1 to 44.1.6.

Security Systems

Popek, G. and D. Farber, "Program Verification Assertions for a Practical Application, Part I," *Proceedings of the Symposium on Distributed Computing Systems: Micros, Minis, and Networks*, Pisa, Italy, December 1975.

Walton, E.J. "The UCLA Pascal Translation System," UCLA-ENG-7954. Also published as Technical Report 76-1, University of California, Los Angeles, January 1976.

Popek, G.J. and D.A. Farber, "On Computer Security Verification," *Twelfth IEEE Computer Society International Conference: Compcon '76*, San Francisco, California, February 1976.

Farber, D. and G. Popek, "An Algebraically Based Programming Language," *Proceedings of the Fourth International Conference on Algorithmic Languages*, Courant Institute, New York, June 1976.

Kemmerer, R. A., "Assignments and Predicates in Kal-Kan," *Fourth International Conference on the Implementation and Design of Algorithmic Languages*, New York, June 1976.

Schwartz, R., "Security Flaws in Tenex, a Time Sharing Operating System," M.S. in Computer Science, Computer Science Department, University of California, Los Angeles, June 1976.

Farber, D.A. "A Model of Program Translation," Computer Science Department, University of California, Los Angeles, Technical Report 76-3, August 1976.

Popek, G.J., J.J. Horning, B.W. Lampson, R.L. London and J.G. Mitchell. "The Programming Language Euclid," Computer Science Department, University of California, Los Angeles, Technical Report 76-4, September 1976.

Walker, B.J. "Note on Proof of Concrete Code," Computer Science Department, University of California, Los Angeles, Technical Report 76-5, September 1976.

Abraham, S.M. "A Protection Design for the UCLA Security Kernel", M.S. in Computer Science, Computer Science Department, University of California, Los Angeles, December 1976. Also published as Technical Report 76-6, UCLA, December 1976.

Popek, G.J., J.J. Horning, B.W. Lampson, J.G. Mitchell and R.L. London. "Notes on the Design of EUCLID," *ACM Conference on Language Design for Reliable Software*, Raleigh, North Carolina, March 1977, pp. 11-18. Also published in *SIGPLAN Notices*, 12(3):11-18, 1977. Also published as Technical Report 77-1, University of California, Los Angeles, March 1977.

Kline, C.S. and G.J. Popek. "Encryption in Computer Network Security," Computer Science Department, University of California, Los Angeles, Technical Report 77-2, April 1977.

Satellite Communications

Nelson, L. and L. Kleinrock, "F-TDMA RFNM Driven Traffic Experiments and Simulations," Packet Satellite Program Working Notes, #18, June 1976.

Gerla, Mario, "S-ALOHA Stability Control Algorithms," Packet Satellite Program Working Notes, #36, September 1976.

Gerla, Mario, "S-ALOHA Stability Control Algorithms and Simulation Experiments," Packet Satellite Working Notes, #46, December 1976.

Nelson, L., "Thoughts on Multi-Net Experiment Control and Data Collection Facilities," Packet Satellite Program Working Note, #56, February 1977.

Gerla, M., L. Nelson, and L. Kleinrock, "F-TDMA, R-TDMA, and S-ALOHA Measurements in SAT-NET," Packet Satellite Program Working Note, #72, May 1977.

Packet Radio Studies

Kleinrock, L. and F. Tobagi, "Packet Switching in Radio Channels: Part I - Carrier Sense Multiple-Access Modes and their Throughput-Delay Characteristics," *IEEE Transactions on Communications*, Vol. COM-23, No. 12, December 1975, pp. 1400-1416.

Tobagi, F. and L. Kleinrock, "Packet Switching in Radio Channels: Part II - the Hidden Terminal Problem in Carrier Sense Multiple-Access and the Busy-Tone Solution," *IEEE Transactions on Communications*, Vol. COM-23, No. 12, December 1975, pp. 1417-1433.

Tobagi, F., "Analysis and Simulation of Buffered Packet Radio Systems," Packet Radio Temporary Note #161, December 1975.

Tobagi, F., "On Measurement Facilities in Packet Radio Systems," Packet Radio Temporary Note #166, January 1976.

Tobagi, F. and L. Kleinrock, "On the Analysis of Buffered Packet Radio Systems," *Proceedings of the Ninth Hawaii International Conference on System Sciences*, Honolulu, Hawaii, January 1976, pp. 42-45.

Tobagi, F. and S. Lieberman, "Measurement in Packet Radio Systems: Methods for Collection of Cumstat Data," Packet Radio Temporary Note #175, March 1976.

Tobagi, F., S. Lieberman and L. Kleinrock, "On Measurement Facilities in Packet Radio Systems," *AFIPS Conference Proceedings*, Vol. 45, National Computer Conference, New York, June 1976, AFIPS Press, Montvale, New Jersey, 1976, pp. 589-596.

Tobagi, F. and L. Kleinrock, "Packet Switching in Radio Channels: Part IV--Stability Considerations and Dynamic Control in Carrier Sense Multiple Access," Packet Radio Temporary Note #181, June 1976.

Tobagi, F. and L. Kleinrock, "Packet Switching in Radio Channels: Part III - Polling and (Dynamic) Split-Channel Reservation Multiple Access," *IEEE Transactions on Communications*, Vol. COM-24, No. 8, August 1976, pp. 832-844.

Tobagi, F. and L. Kleinrock, "Packet Switching in Radio Channels: Effect of Error Control Traffic on ALOHA and CSMA Channel Performance," Packet Radio Temporary Note #201, November 1976.

Tobagi, F., "NETCAP: A Simulator for PR Net Channel Access Protocol Performance," Packet Radio Temporary Note, #203, December 1976.

Tobagi, F. and Y. Yemini, "NETCAP II: Derivation of Confidence Intervals for Various Performance Measures," Packet Radio Temporary Note, #204, December 1976.

Scholl, M. and L. Kleinrock, "On A Mixed Mode Multiple Access Scheme for Packet-Switched Radio Channels," Packet Radio Temporary Note #214, March 1977.

Kleinrock, L. and M. Scholl, "Packet Switching in Radio Channels- New Conflict-Free Multiple Access Schemes," Packet Radio Temporary Note #225, June 1977.

Kleinrock, L. and M. Scholl, "Packet Switching in Radio Channels: New Conflict-Free Multiple Access Schemes for a Small Number of Data Users," *Conference Record, International Conference on Communications*, Vol. 2, Chicago, Illinois, June 1977, pp. 22.1-105 to 111. (also issued as Packet Radio Temporary Note #210, January 1977.)

M. Scholl and L. Kleinrock, "On The M/G/1 Queue with Rest Period and Certain Service-Independent Queueing Disciplines," Packet Radio Temporary Note #226, June 1977.

Dissertations

Kamoun, F., "Design Considerations for Large Computer Communication Networks," Ph.D. Dissertation, Computer Science Department, University of California, Los Angeles, March 1976.

Scholl, M., "Multiplexing Techniques for Data Transmission over Packet Switched Radio Systems." Ph.D. Dissertation, Computer Science Department, University of California, Los Angeles, September 1976.

Erlich, Z., "On Central Bus Transportation Systems with Poisson Arrivals," Ph.D. Dissertation, Computer Science Department, University of California, Los Angeles, September 1976.

Naylor, W., "Stream Traffic Communication in Packet Switched Networks," Ph.D. Dissertation, Computer Science Department, University of California, Los Angeles, September 1977.

Kermani, P., "Switching and Flow Control Techniques in Computer Communication Networks," Ph.D. Dissertation, Computer Science Department, University of California, Los Angeles, December 1977.

Principles and Lessons in Packet Communications

LEONARD KLEINROCK, FELLOW, IEEE

Invited Paper

Abstract—After nearly a decade of experience, we reflect on the principles and lessons which have emerged in the field of packet communications. We begin by identifying the need for efficient resource sharing and review the original and recurring difficulties we had in achieving this goal in packet networks. We then discuss various lessons learned in the areas of: deadlocks; degradations; distributed control; broadcast channels; and hierarchical design. The principles which we discuss have to do with: the efficiency of large system; the switching computer; network constraints; distributed control; flow control; stale protocols; and designers not yet experienced in packet communications. Throughout the paper, we identify various open issues which remain to be solved in packet communications.

I. INTRODUCTION

WHAT IS IT WE now know about packet communications that we did not know a decade ago? What made the problem difficult, and why were the solutions not immediately apparent to us in the late 1960's? Whereas the answers to these questions may entice the system designer (indeed, I, for one, delight in such investigations), why should the network user care a whit? To most users (and, alas, to many designers), communications is simply a nuisance and they would just as soon ignore those problems and get on with the "real" issues of information processing!

In this paper (and in conjunction with the other papers in this Special Issue of the PROCEEDINGS), we hope to answer some of these questions. We will identify the need for resource sharing, explain why the problem of efficient resource sharing is hard, and why it must be understood, review some of the lessons we learned (mostly from the three ARPA packet networks), and then, finally, state some principles which have evolved from the study and extensive use of packet communications.

II. RESOURCE SHARING

A privately owned automobile is usually a waste of money! Perhaps 90 percent of the time it is idly parked and not in use. However, its "convenience" is so seductive that few can resist the temptation to own one. When the price of such a poorly utilized device is astronomically high, we do refuse the temptation (how many of us own private jet aircraft?). On the other hand, when the cost is extremely low, we are obliged to own such resources (we all own idle pencils).

An information processing system consists of many poorly utilized resources. (A resource is simply a device which can perform work for us at a finite rate.) For example, in an information processing system, there is the CPU, the main memory, the disk, the data channels, the terminals, the printer etc.

Manuscript received March 24, 1978; revised June 16, 1978. This research was supported by the Advanced Research Projects Agency of the Department of Defense under Contract MDA 903-77-C-0272.

The author is with the Computer Science Department, University of California, Los Angeles, CA 90024.

One of the major system advances of the early 1960's was the development of multiaccess time-sharing systems in which computer system resources were made available to a large population of users, each of whom had relatively small demands (i.e., the ratio of their peak demands to their average demands was very high) but who collectively presented a total demand profile which was relatively smooth and of medium to high utilization. This was an example of the advantages to be gained through the smoothing effect of a large population (i.e., the "law of large numbers") [1]. The need for resource sharing is present in many many systems (e.g., the shared use of public jet aircraft among a large population of users).

In computer communication systems we have a great need for sharing expensive resources among a collection of high peak-to-average (i.e., "bursty") users [1]. In Fig. 1 we display the structure of a computer network in which we can identify three kinds of resources:

1) the terminals directly available to the user and the communications resources required to connect those terminals to their HOST computers or directly into the network (via TIPS in the ARPANET, for example—this is an expensive portion of the system and it is generally difficult to employ extensive resource sharing here due to the relative sparseness of the data sources;

2) the HOST machines themselves which provide the information processing services—here multiaccess time sharing provides the mechanism for efficient resource sharing;

3) the communications subnetwork, consisting of communication trunks and software switches, whose function it is to provide the data communication service for the exchange of data and control among the other devices.

The HOST machines in 2) above contain hardware and software resources (in the form of application programs and data files) whose sharing comes under the topic of time sharing; we dwell no further on these resources. Rather, we shall focus attention on those portions of the computer communications system where packet communications has had an important impact. Perhaps the most visible component is that of the communications subnetwork listed in item 3) above. Here packet communications first demonstrated its enormous efficiencies in the form of the ARPANET in the early 1970's (the decade of computer networks). The communication resources to be shared in this case are *storage capacity* at the nodal switches (these switches are called IMP's in the ARPANET), *processing capacity* in the nodal switches, and *communications capacity* of the trunks connecting these switches. Packet switching in this environment has proven to be a major technological breakthrough in providing cost effective data communications among information processing systems. Deep in the backbone of such packet-switched networks there is a need for long-haul high-capacity inexpensive communications, and it is here where we

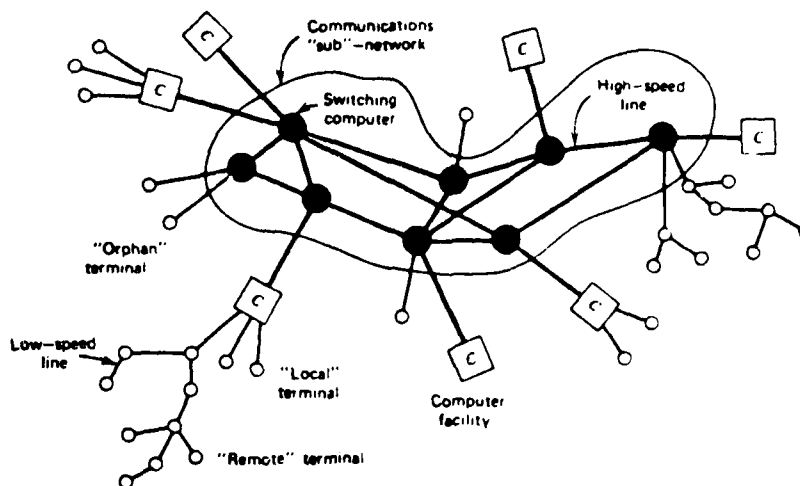


Fig. 1. The structure of a computer-communication network.

see the second application of packet communications for resource sharing in the form of satellite packet switching; elsewhere in this PROCEEDINGS [2] you will find a description of the SATNET, an ARPA-sponsored research network connected to the ARPANET. The third application may be found in the local access problem stated in item 1) above which also lends itself to the use of packet switching to provide efficient communications resource sharing; this takes the form of the use of a multiaccess broadcast channel in a local environment, commonly known as ground radio packet switching. Here too, ARPA has sponsored an experimental system, and its description may be found in this PROCEEDINGS [3]. The common element running through all these systems is the application of the smoothing effect of a large population to provide efficient resource sharing, an exquisite example of which is provided by packet communications.

Having described the environment and the resources of interest, let us now discuss the performance measures which permit us to evaluate the effort of resource sharing in a quantitative way. Indeed, there are basically four measures that both the system designer and network user apply in evaluating the service provided in a communications environment. These are *throughput*, *response time*, *reliability*, and *cost*. Before packet networks came into existence, the obvious solution for providing communications between two devices was to lease or dial a line between the two. In such a case the user was able to associate precise quantitative values to the four measures listed above. On the other hand when one attaches to a packet network, the user cannot get deterministic answers to the same quantities as he has in the past. He must accept probabilistic statements regarding throughput, delay, and reliability (and alas, sometimes even cost). Moreover the quantities so prescribed can seldom be measured in a straightforward fashion. This is the state of affairs to which we have evolved today! It is to the credit of those who developed packet communications in the last decade that the system design was carefully studied and well-analyzed *prior* to and during the systems implementation; this certainly has not, in general, been the history in the information processing industry.

III. WHY THE PROBLEM IS HARD

Back in 1967, when the concept of the ARPANET was first taking form, we found ourselves entering the uncharted terrain

of packet switching. Let us trace our initial confusion regarding that project briefly. Certainly, there existed at that time some communication networks, but they were mostly highly specialized networks with restricted goals. In the early 1960's Paul Baran had described some of the properties of data networks in a series of Rand Corporation papers [4]. He focused on the routing procedures and on the survivability of distributed communication systems in a hostile environment, but did not concentrate on the need for resource sharing in its form as we now understand it; indeed, the concept of a software switch was not present in his work. In 1968 Donald Davies at the National Physical Laboratories in England was beginning to write about packet-switched networks [5]; at around the same time, Larry Roberts at ARPA pursued the use of packet switching in an experimental nationwide network [6]. For a more complete history of the evolution of packet communications, see [7].

In the initial conception of a packet network, we identified some problems and looked to the technical literature for solutions to these problems. For example, how should one design the topology of a network, and how should one select the bandwidth for the various channels in such a network, and in what fashion should one route the data through the network, and what rules of procedure should two communicating processes adopt, and how much storage did one need at the multiplexing nodes of the network? These and many other questions confronted us. Indeed the general problem was how to achieve efficient resource sharing among a set of incompatible devices in a geographically distributed environment where access to these devices arose from asynchronous processes in a highly bursty fashion. Moreover, not only was the demand process bursty, it was also highly unpredictable in the sense that the instants when the demands arose and the duration of the demands were unknown ahead of time. Fortunately we were unaware of the enormity of the problems facing us and so we plunged ahead enthusiastically and with naive optimism. The remainder of this section describes why the problem was difficult, and in following sections we describe the lessons we learned and the principles we established in the development of packet communications. Our efforts have been well rewarded and the technology of packet communications has come of age and has proven itself to be a cost-effective technology.

We quickly found that many of our old techniques could not be directly applied to packet network design and that new techniques had to be developed; these new techniques turned out to be of great generality and have led to principles and to understanding which are sure to benefit us for many years to come. One of our first tasks was to develop tools which would allow us to analyze the performance of a given network. This involved evaluating the delay-throughput profile for networks. Basically, this is a queueing problem in a network environment. It had earlier been recognized [8] that the probabilistic complexities one encounters in computer networks are extremely difficult. One of the simplest analytical questions involving the solution of two nodes in tandem was first posed at that time in 1964 and has only been satisfactorily answered (in the queueing theoretic sense) within the past year [9]; this, note, is for the simplest problem. Indeed we have come to realize that an exact solution for the delay-throughput profile is probably hopeless in a complex network environment. Fortunately suitable approximations [1], [8] have been developed which permit one to predict the performance of given networks with a high degree of accuracy. More than that, these approximate tools allow us to expose and understand the phenomenological behavior of networks.

The astute reader will observe that the resource sharing problem stated above sounds very much like the problem faced in the design of time sharing systems. Surely, with time sharing, we are faced with the problem of sharing resources among asynchronous processes which behave in a bursty fashion. The major difference between the two problems, however, is that our problem exists in a *geographically distributed* environment which requires expensive communication resources in the communications and coordination functions. The implications here are strong. For example, when communication is cheap, then wide-band communications can be obtained with extremely small delays; such is the case, for example, within the resources of a local operating system connected together by a data bus. In a regional or nationwide network subject to the relatively expensive cost of telecommunications, we find that typical bandwidths are many orders of magnitude less than that in a local time-sharing environment, and the delays are many orders of magnitude greater. Furthermore, the control of these processes in the time sharing environment can be very tightly coupled if desired or left loosely coupled if there is sufficient reason; in the network environment we typically find our processes are very loosely coupled due to the difficulty of tightening the control between them (indeed, the inherent delay due to the finite speed of light is a fundamental limitation on the tight coupling of remote processes). The overhead in the time sharing system is variable and may be very high with poor system design (for example, thrashing) but may be made small with clever design. In the network environment, for a variety of reasons, we find that the overhead due to packet headers, control information and resource allocation tends to be relatively high. These comparisons are summarized in Table I.

Not only do we have extremes in communications cost between these two systems, we also have a significant difference in the reliability of that communications. Indeed, in the local time shared system, the process-to-process communication is usually assumed to be reliable and therefore the acknowledgment procedure (if any exists) is simple and tends to be invoked only under exceptional circumstances. On the other hand, in the distributed computer network environment, communications reliability is not assumed, and, therefore, an elaborate

TABLE I
ASYNCHRONOUS PROCESS-TO-PROCESS COMMUNICATION AND CONTROL
COST COMPARISON BETWEEN LOCAL PROCESSES IN A TIME-SHARED
SYSTEM AND DISTRIBUTED PROCESSES IN A NETWORK

	Multiaxess Time-Shared Systems	Geographically Distributed Computer Networks
Typical bandwidth	megabytes/sec	kilobits/sec
Typical communications delay	fractions of a micro-second	tens to hundreds of milliseconds
Process-process coupling	tight to loose	typically loose
Overhead due to system control	variable (typically low)	variable (typically high)

acknowledgment procedure (often including timeouts and other defaults) is usually invoked. We are here reminded of the "two-army" problem. This is the problem where two blue armies are each poised on opposite hills preparing to attack a single red army in the valley. The red army can conquer either of the blue armies separately but will fall to defeat if both blue armies attack simultaneously. The blue armies communicate with each other over an unreliable communications system. The problem is to coordinate the two blue armies so that they will attack simultaneously. Let us assume that Blue Army 1 (B1) sends a message to Blue Army 2 (B2) indicating that they should jointly attack at noon the next day. Clearly B1 must await a positive acknowledgment from B2 that the command was properly received; were B1 to attack without such an acknowledgment, then the possibility exists that B2 did not receive the message correctly, in which case B1 is subject to the certain annihilation of his army. If B2 properly receives the command and acknowledges it, then he must await an acknowledgment of the acknowledgment from B1, for if B1 did not receive his acknowledgment then we know B1 will not attack and B2's attack will then be doomed to failure. The argument continues in a circular fashion where acks of acks . . . are continually transmitted with no action ever being taken; the two blue armies can *never* get perfectly synchronized with certainty using this unreliable communications system.

We see therefore that the new culprit in resource sharing in a distributed environment is the fact that the resources are distributed and cannot easily be assigned to the demands. Indeed, in previous resource allocation problems, which often come under the name of scheduling algorithms, we made a big assumption, namely, that the scheduling information could be passed around for free. That is, the competing demands could organize themselves into a cooperating queue at no price. Unfortunately, in the distributed environment, the cost of organizing demands into a cooperating queue may be very large, and in one fashion or another nature will make you pay a price [10].

The problem of resource sharing in a distributed environment manifests itself in the routing control and flow control problems. The problem of flow control is to regulate the rate at which data crosses the boundary of the communications sub-network (both into and out of the network). The problem of routing control is to efficiently transport the data (which has been admitted by the flow control procedure) across the network to its destination. In 1967 we were aware of the sophistication needed in the routing procedure, but were relatively ignorant of the need for effective flow control (see below). These two problems are difficult because we are dealing with a

control procedure in a distributed environment subject to random delays in passing that control information around in order to control the random demands. The purpose of both procedures is to efficiently use the network resources (IMP storage, IMP processing capacity, and communications capacity). In achieving this goal one must attempt to control congestion, route data around busy or defective portions of the network, and in general must adaptively assign capacity to the data traffic flow in an efficient, dynamic way. These are hard control problems and represent a class which has not been adequately studied up to and including the present time. We have come to learn that distributed control is a sophisticated problem. Below we return to the issues involved in distributed resource allocation and sharing. For the moment let us introduce some of the other sources of complexity in packet communications.

In any distributed communications system design one is faced with a *topological design problem*. The problem basically is, given a set of constraints to meet, find that topological design structure which meets these constraints at least cost. The field of network flow theory addresses itself to such problems and the salient feature of this theory is that most of its problems cannot be solved! To exhaustively search over all possible topological designs for a given problem is certainly not a solution since the number of possible alternatives to consider can easily exceed the number of atoms in the universe even for relatively small problems. (For example, if at some stage you must consider all permutations of 20 objects, then a computer would take more than 75 000 years to process all $20!$ cases even if it could examine one million cases per second.) Rather, a solution consists of elegant search procedures which are computationally efficient and which find the optimal topology for the given problem. Very few problems in network flow theory yield to such efficient algorithms. Rather, one gets around the combinatorial complexity naturally inherent in these problems by accepting suboptimal solutions. (Beware! A suboptimal solution to a problem is simply the result of a procedure which examines a subset of possible solutions and picks the best of those examined—this suboptimal solution may or may not be close to the optimal.) The trick here is to find efficient heuristic search procedures which come close to the optimal rapidly. Over the past decade, efficient procedures have been developed in many cases and new procedures are constantly being investigated for the topological design problem.

Another source of difficulty in the resource sharing problem is in defining the appropriate *performance measure*. For example, what is the capacity of a network? It is well-known in network flow theory that one can easily calculate the capacity (i.e., the throughput) between any two pairs of points. What is not straightforward is to evaluate the total data-carrying capacity of a network where throughput is measured in terms of messages successfully received at their destination. The difficulty comes about because the capacity of the network strongly depends upon the traffic matrix one assumes for the data flowing through that network. For example, if the traffic matrix were such that traffic passes only between immediate neighbors in the topological structure and in an amount equal to the capacity of the line connecting those two, then the network capacity would approach a value equal to the sum of all channel capacities in the network. This is clearly an upper bound to the capacity for any other traffic matrix. Since in general we do not know the traffic matrix for a network to be designed for future use, how is one to evaluate that capacity?

Yet another source of difficulty in the network problem is that of *interfacing terminals and HOSTS to networks as well as one network to another network*. For example, there is the general issue of *virtual circuit service* as opposed to *datagram service*. A virtual circuit service presents to the user a communications environment in which all functions appear as if she were directly connected between the two points communicating (i.e., an orderly and controlled flow is maintained whereby data is guaranteed to be correct upon delivery to the destination, comes out of the network in the same order in which it came in, and a flow control may be applied to that transmission to maintain efficient use of network resources and of terminal-HOST resources). A datagram service ensures none of these things and simply sends blocks of data (packets) across the network, not guaranteeing correct delivery (or delivery at all), not guaranteeing orderly flow in terms of sequencing (packets may arrive at the destination out of order), and not enforcing any flow control procedure at the process-to-process level. Which of those two services is desirable has become an issue of international proportions discussed elsewhere in these proceedings [7], [11]. Furthermore, the interconnection of two networks presents an enormous and rich variety of difficult problems. For example, should one apply flow control at the boundary of each network in a tandem chain of networks or should one apply flow control from the source HOST to the destination HOST across many networks simultaneously? If we consider the interconnection of networks with different packet sizes, we have the general problem of fragmentation whereby long packets get fragmented into smaller packets when crossing network boundaries. A variety of other very important issues arise in interconnecting (see [27] for a more detailed discussion of interconnecting).

Many of the problems we have just presented come about due to the distributed nature of our message sources and system resources. The problems created by distributed sources are very clearly seen in the environment of geographically dispersed message sources which communicate with each other over a common broadcast channel; in this case, access to the capacity of the channel must be carefully coordinated.

Thus in answering the question, "Why is the problem hard?" we have found the following sources of difficulty:

- 1) the analytic problems from queueing theory and the probabilistic complexity resulting therefrom;
- 2) the topological design problems from network flow theory and the combinatorial complexity resulting therefrom;
- 3) the price for coordinating and sharing resources in a geographically distributed environment (the new culprit) leading to problems of resource allocation, routing control, flow control, and general process-to-process communication problems.

IV. LESSONS LEARNED

After a decade of experience with packet communications it is fair that we ask what lessons have we learned and what have we come to know about the needs of the user and the questions he would like to have answered. So far as the user is concerned we shall see as we step through the lessons learned below that he cannot insulate himself completely from the underlying technology of packet communications. Indeed the service he sees is quite different from that which he has with leased lines as mentioned above. Moreover, certain decisions will either be thrust upon him or accepted by him due to the nature of the service offered; if he is unaware of the consequence of setting

parameters in that decision-making process then he may seriously degrade the performance of the network due to his ignorance. Let us now list some of the lessons we learned and return to the principles in the following section.

A Deadlocks

In [1], [12], and [13] we described in detail some of the deadlocks and degradations of which we have become aware. In this section we simply enumerate and sketch the details of the deadlocks. Simply stated, a *deadlock* (also commonly referred to as a *lockup*) is the unpleasant situation in which two (or more) competing demands have each been assigned a subset of their necessary resources; neither can proceed until one of them collects some additional resources which currently are assigned to the other and neither demand is willing to release any resource currently assigned to him. Deadlocks are one of the most serious system malfunctions possible, and one must take great care to avoid them or find ways to recover from them. It is ironic that flow control procedures by their very nature present *constraints* on the flow of data (e.g., the requirement for proper sequencing), and if the situation ever arises whereby the constraint cannot be met, then, by definition, the flow will stop, and we will have a deadlock! This is the philosophical reason why flow control procedures have a natural tendency to introduce deadlocks. In this section we briefly discuss four ARPANET deadlocks which have come to be known as: *reassembly lockup*; *store-and-forward deadlock*; *Christmas lockup*; and *piggyback lockup*.

Reassembly lockup, the best known of the ARPANET deadlock conditions (and one which was known to exist in the very early days of the ARPANET implementation), was due to a logical flaw in the original flow-control procedure. In the ARPANET, a string of bits to be passed through the network is broken into "messages" which are at most approximately 8000 bits in length. These messages are themselves broken into packets which are at most approximately 1000 bits in length. A message requiring more than one packet (up to a maximum of eight) is termed a multipacket message and each of these packets traverses the network independently; upon receipt at the destination node, these packets are "reassembled" into their original order and the message itself is recomposed, at which time it is ready for delivery out of the network. (A more complete description of the ARPANET protocols may be found in [1], [13].) Reassembly lockup occurred when partially reassembled messages could not be completely reassembled since the network through which the remaining packets had to traverse was congested, thus preventing these packets from reaching the destination; that is, each of the destination's neighbors had given all of their relay (store-and-forward) buffers to additional packets (from messages other than those being reassembled) heading for that same destination and for which there were no unassigned reassembly buffers available. Thus the destination was surrounded by a barrier of blocked IMP's which themselves could provide no store-and-forward buffers for the needed outstanding packets, and which at the same time were prevented from releasing any of their store-and-forward buffers since the destination itself refused to accept these packets due to a lack of reassembly buffers at the destination. The deadlock was simply that the remaining packets could not reach the destination and complete the reassembly until some store-and-forward buffers became free,

and the store-and-forward buffers could not be released until the remaining packets reached the destination.

Store-and-forward deadlock is another example of a lockup that can occur in a packet-switched network if no proper precautions are taken [1], [13]. The case of "direct" store-and-forward lockup is simply a "stand-off." Let us assume that all store-and-forward buffers in some IMP *A* are filled with packets headed toward some destination IMP *C* through a neighboring IMP *B* and that all store-and-forward buffers in IMP *B* are filled with packets headed toward some destination IMP *D* through IMP *A*. Since there is no store-and-forward buffer space available in either IMP *A* or *B*, no packet can be successfully transmitted between these two IMP's and a deadlock situation results. One way to prevent the deadlock is to prohibit these packets in IMP *A* from occupying all those store-and-forward buffers which are needed by the packets coming in from IMP *B* (and vice versa) by the introduction of "buffer classes" as in [14]. This is accomplished by partitioning the buffers in a switch into classes, say, B_0, B_1, \dots, B_k , where k is the longest path length in the network. A packet arriving at a switch from outside the net has access only to class B_0 buffers. When a packet arrives at a switch after having made k hops so far, it has access to class B_0, \dots, B_k buffers, etc. Thus, the closer a packet gets to its final destination, the more access it has, and therefore the speedier its passage through the network. It can be proven [14] that this "buffer class" allocation will prevent direct store-and-forward lockup. "Indirect" store-and-forward lockup can occur when all store-and-forward buffers in a loop of IMP's become filled with packets all of which travel in the same direction (clockwise or counter-clockwise) and none of which are within one hop of their destination. Both store-and-forward lockup conditions are far less likely if, as in the ARPANET, more than one path exists between all pairs of communicating IMP's.

In December 1973, the dormant *Christmas lockup* condition was brought to life. This lockup was exposed by collecting measurement messages at UCLA from all IMP's simultaneously. The Christmas lockup occurred when these measurement messages arrived at the UCLA IMP for which reassembly storage had been allocated but for which no reassembly blocks had been given. (A reassembly block is a piece of storage used in the actual process of reassembling packets back to messages.) These messages had no way to locate their allocated buffers since the pointer to an allocated buffer is part of the reassembly block; as a consequence, allocated buffers could never be used and could never be freed. The difficulty was caused by the system first allocating buffers before it was assured that a reassembly block was available. To avoid this kind of lockup, reassembly blocks are now allocated along with the reassembly buffers for each multipacket message in the ARPANET.

Piggyback lockup is a deadlock condition which was identified by examining the flow control code and has, as far as we know, never occurred. This lockup condition comes about due to an unfortunate combination of intuitively reasonable goals implemented in the flow-control procedure. One of these goals, which we have already mentioned, is to deliver messages to a destination in the same order that the source received them. The other innocent condition has to do with the reservation of reassembly storage space at the destination. In order to make this reservation procedure efficient, it is reasonable that only the first multipacket message of a long file transfer be required

to make the reservation. The ARPANET flow control procedure will then maintain that reservation for a given file transfer as long as successive multipacket messages from that file are promptly received in succession at the source IMP. We have now laid the groundwork for piggyback lockup. Assume that there is a maximum of eight reassembly buffers in each IMP; the choice of eight is for simplicity, but the argument works for any value. Let IMP A continually transmit eight-packet messages (from some long file) to some destination IMP B such that all eight reassembly buffers in IMP B are used up by this transmission of multipacket messages. If now, in the stream of eight-packet messages, IMP A sends a single-packet message (not part of the file transfer) to destination IMP B it will generally not be accepted since there is no reassembly buffer space available. The single packet message will therefore be treated as a request for buffer allocation (these requests are the mechanism by which reservations are made). This request will not be serviced before the RFNM (an end-to-end acknowledgment from the destination to source) for the previous multipacket message has been sent. When the RFNM is generated, however, all the free reassembly buffers will immediately be allocated to the next multipacket message in the file transfer for efficient transmission as mentioned above; this allocation is said to be "piggybacked" on the RFNM. In this case, the eight-packet message from IMP A that arrives later at IMP B (and which is stored in the eight buffers) cannot be delivered to its destination HOST because it is out of order. The single-packet message that should be delivered next, however, will never reach the destination IMP since there is no reassembly space available, and, therefore, its requested reservation can never be serviced. Deadlock! A minor modification removes the piggyback lockup.

These various deadlock conditions are usually quite easy to prevent once they are detected and understood. The trick, however, is to expurgate all deadlocks from the control mechanism ahead of time, either by careful programming (a difficult task) or by some automatic checking procedure (which may be as difficult as proving the correctness of programs). Those deadlocks found in the ARPANET have, to the best of our knowledge, been eliminated.

B. Degradations

A degradation is just that, namely, a reduction in the network's level of performance. (Deadlocks are, of course, an extreme form of degradation which is why we discussed them in the separate section above.) For our purposes, we shall measure performance in terms of delay and throughput. In this section we discuss four sources of ARPANET degradation, namely: *looping* in the routing procedure; *gaps* in transmission; *single-packet turbulence*; and *phasing*.

Looping comes about due to independent routing decisions made by separate nodes which cause traffic to return to a previously visited node (or, in a more general definition, causes traffic to make unnecessarily long excursions on the way to its destination). Of course any reasonable adaptive routing procedure will detect these loops (through the build-up of queues and delays perhaps) and will then break the loop and guide the traffic directly on to its destination. However, the occurrence of loops does cause occasional large delays in the traffic flow and in some applications this is quite unacceptable. It is ironic that a remedy which was introduced in the ARPANET to reduce the occurrence of loops, in fact made them worse in

the sense that whereas they occurred less frequently, when they did occur, they persisted for a longer time. Some loop-free algorithms have recently been published [15], [16].

The next degradation we wish to discuss is the occurrence of *gaps* in the message flow. These gaps come about due to a limitation on the number of messages in transit which the network will allow. Assume that between any source and destination, the network will permit n messages in flight at a time. If n messages are in flight, then the next one may not proceed until an end-to-end acknowledgment is returned back at the source for any one of the n outstanding messages. We now observe that if the round-trip delay (i.e., the time required to send a message across the network in the forward direction and to return its acknowledgment in the reverse direction) is greater than the time it takes to feed the n messages into the network, then the source will be blocked awaiting ack's to release further messages. This clearly will introduce gaps in the message flow resulting in a reduced throughput which we might classify as a mild form of degradation.

We now come to the issue of *single packet turbulence* as observed in the ARPANET. We note that "regular" single-packet messages in the early ARPANET were not accepted by the destination if they arrived out of order. Rather, they were then treated as a request for the allocation of one reassembly buffer. Therefore if, in a stream of single-packet messages, packet p arrived out-of-order (say it arrived after packet $p + 3$), then packets $p + 1$, $p + 2$, and $p + 3$ would all be discarded at the destination, and only after packet p arrived would a single packet buffer be allocated to message $p + 1$. This allocation piggybacked on the end-to-end ack for packet p , and when it arrived at the source IMP, it then caused a retransmission of the discarded packet $p + 1$ (which had been stored in the source). Of course any packets arriving at the destination after packet $p + 3$, but before $p + 1$ arrived in order, would themselves be discarded. When packet $p + 1$ finally arrived for the second time at the destination IMP it was then in order and this caused an allocation of a single-packet buffer to packet $p + 2$, etc. The net result was that only one packet would be deliverable to the destination per round-trip time along this path; had no packets been received out-of-order, then we would have been pumping at a rate close to n packets per round trip time (if the maximum number in transit n could fit into the pipe). Observe that once a single packet arrived out-of-order in this stream, then the degradation from n to 1 packets per round-trip time would persist forever until either some supervisory action was taken or until the traffic stream ceased and began again from a fresh start in the future. We refer to this effect as "single-packet turbulence," and it was observed in the ARPANET as described in [17]. The need to handle a continuous stream of traffic (e.g., packetized speech) was recognized some time ago and resulted in the definition of "irregular" packets known as type 3 packets (as contrasted to "regular" type 0 packets); these packets are allowed to be delivered out of order, receive no end-to-end acknowledgment, and are generally handled in a much more relaxed fashion.

The last degradation we discuss is known as "*phasing*." In a typical packet network, more than one resource is often required before a message is allowed to flow across that network. For example, some required resources may be: a message number; storage space at the source; storage space at the destination, etc. Tokens move around the network passing out

these resources in some distributed fashion. Phasing is the phenomenon whereby enough free tokens are available in the network to permit message flow, but, the proper mix of tokens is not available simultaneously at the proper location in the net. The delay in gathering these tokens represents a degradation [1], [18].

Fortunately, the degradations here described have been remedied in the ARPANET and in later networks.

C. Lessons of Distributed Control

We have had "lessons" in two areas of distributed control. The first has to do with flow control, and it is simply the observation that flow control procedures are rather difficult to invent and extremely difficult to analyze. The deadlocks and degradations referred to in the previous subsections were principally due to flow control failures (and occasionally routing control failures). To date there is no satisfactory theory or procedure for designing efficient flow control procedures, much less evaluating their performance, proving they contain no deadlocks, and proving that they are correct. Attempts in this direction are currently under way.

An important lesson we have learned with flow control is that a packet communications system offers an opportunity for passing data between two devices of (very) different speeds. We can effectively connect a slow-speed teletype to an enormously high-speed memory channel over a packet network and apply flow control procedures which protect the two devices from each other as well as protecting the net from both. Specifically, we must not drown the teletype with a flood of high-speed input, nor must we "nickel-and-dime" a high performance HOST to death with incessant interrupts, nor must we use the network as a storage medium for megabytes of data. Flow control mechanisms provide the means to accomplish this; the trick is to do it well.

The second area of distributed control in packet communications has to do with the routing control. The ARPANET, and many of the networks which have since based their design on the packet-switching technology which emerged from the ARPANET experiment, employ an adaptive routing procedure with distributed control. In such a procedure, routes for the data traffic are not preassigned but rather are dynamically assigned when they are needed according to the current network status. Control packets (called routing update packets) which describe the state of the network to some degree are passed back and forth between neighboring IMP's in some fashion and current queue lengths and congestion measures are added to these updates by each IMP. The ARPANET employs a periodic update routing procedure whose rate depends upon channel utilization and line speed. The updates passing between IMP's have no priority in competing for the processing capacity of the CPU at the IMP's but do have priority in the queue discipline feeding the output modems between IMP's. An important lesson learned is that giving low priority to the processing of routing updates appears to be advantageous since the processing load on the CPU is rather large and prevents the further dispatching of arriving packets to output queues [19]. Another routing lesson we have learned is that frequent updates cause background congestion in a network which may be intolerably high even in the absence of other data traffic; the update procedure and update rate must be carefully chosen. A number of alternatives to periodic updates have been suggested [1] including such things as aperiodic updating (send updates only when status information has crossed certain thresholds and then send it immediately); and purely local information for

routing decisions based on queue lengths within a given node and knowledge of the current topology. Furthermore, unless care is taken, there is a tendency for looping to occur in these distributed control algorithms; looping can be prevented with more sophisticated algorithms [15].

One of the lessons which is now beginning to emerge is that the most important advantage of distributed control adaptive routing is its ability to automatically sense configuration changes in the network; these configuration changes may be planned or accidental as for example the result of a line or IMP failure. This is important for two reasons: first because configuration changes do happen often enough so that the requirement for a centralized control evaluating new routing tables based on the current configuration would be an enormously complex task from an administrative point of view; second because it is specifically at times of configuration changes when drastic network action must be taken and only then is the adaptive routing procedure really called upon to do serious work (it is not yet clear to what extent the routing algorithm should adapt to statistical fluctuations in traffic).

Without diminishing the result of these lessons, it is fair to say that the most significant lesson learned regarding routing is that it works at all. Perhaps one of the greatest successes of the ARPANET experiment was to show that a distributed control adaptive routing algorithm would indeed converge on routes which were sufficiently good. The difficulty in proving this lies in the fact that we are dealing with a dynamic situation in a distributed control environment with delays in the feedback paths for control information flow. The empirical proof that things do work has had an important impact on network design; indeed, these distributed algorithms are currently operating successfully in a number of packet networks.

D. Lessons from Broadcast Channels

As mentioned earlier, packet communications has found important applications in the areas of satellite packet broadcasting and in ground radio packet switching. In both environments we have a situation in which a common broadcast channel is available to be shared by a multiplicity of users. Since these users demand access to the channel at unpredictable times, we must introduce some access scheme to coordinate their use of the channel in a way which prevents degradations and mutual interference. In many of the schemes described [10] we have found that "burst" communications provides efficiencies over that of "trickle" transmission. By this we mean that when a data source requires access to the channel, it should be given access to the full capacity of that broadcast channel and not be required to transmit at a slow speed using only a fraction of the available bandwidth (see Section V-A on principles regarding "bigger is better").

In examining the recent literature, we find that a number of access schemes have been invented, analyzed, and published; for a summary of many of these access schemes, see [10]. We observe that these access schemes fall into one of three categories, each with its own cost. The first of these involves random access contention schemes whereby little or no control is exerted on the users in accessing the channel, and this results in the occasional collision of more than one packet; a collision destroys the use of the channel for that transmission. Such schemes as pure ALOHA, slotted ALOHA, and (to a much lesser extent) Carrier Sense Multiple Access fall into this category. At the opposite extreme, we have the static reservation access methods which preassign capacity to users thereby

creating "dedicated" as opposed to multiaccess channels. Here the problem is that a bursty user will often not use his preassigned capacity in which case it is wasted. Such schemes as Time Division Multiple Access and Frequency Division Multiple Access fall in this category. Between these two extremes are the dynamic reservation systems which only assign capacity to a user when he has data to send. The loss here is due to the overhead of implementing the demand access. Such schemes as Polling (where one waits to be asked if he has data to send), active reservation schemes (where one asks for capacity when he needs it), and Mini-Slotted Alternating Priority (where a token is passed among numbered users in a prearranged sequence, giving each permission to transmit as he receives a token) all fall in this category. Each of these schemes pays its tribute to nature as shown in Table II.

Unfortunately, at this point in time we are unable to evaluate the minimum price (i.e., a degradation to throughput and/or delay) one must pay for a given distributed multiaccess broadcast environment.

We have found that contention schemes are fundamentally unstable in that they have a tendency to drift into a congested state where the throughput decreases significantly at the same time the delay increases. Fortunately, however, we have been able to design and implement amazingly effective control schemes which stabilize these contention schemes [20]. Another lesson we have learned is that certain tempting ways of mixing two access schemes (e.g., taking a fraction of the traffic and a fraction of the capacity assigned to one access scheme, and using that capacity to handle that traffic using a second access scheme) does not give an improvement in the overall throughput-delay performance [10]. Furthermore we have found that certain capture effects exist in some of the contention schemes (e.g., a group of terminals may temporarily hog the system capacity and thereby "lock out" other groups for extended periods of time) and one must be wary of such phenomena [20].

We have also found that in a ground radio broadcast environment, a few buffers in each packet radio unit appear to be sufficient to handle the storage requirements [21]; this comes about largely due to the fact that our transceivers are half-duplex (i.e., they can either transmit or receive, but not both, at a given time). We can show (see Section IV-E) that dedicated broadcast channels have an inherent advantage over dedicated wire networks in a large (many-user) bursty store-and-forward environment [22]. Moreover, we have investigated the optimal transmission range for ALOHA networks and have found that those broadcast networks can be made quite effective when the traffic is not bursty; indeed this optimal range is chosen so that the channel utilization in the resulting local ALOHA system is $1/2e$ and then those networks need only \sqrt{e} more capacity than the corresponding M/M/1 network [22].

Lastly we point out that perhaps one of the first applications of broadcast radio access schemes will be to implement these access schemes on wire networks (for example, coaxial cables or fiber-optics channels) in a local environment; an example of such an implementation is the Ethernet [23].

E. Hierarchical Design

As N (the number of nodes in a network) grows, the cost of creating the topological design of such a network behaves like N^E where E is typically in the range from 3 to 6. Thus we see that topological design quickly becomes unmanageable. Secondly, we note that as N grows, the size of the routing table in each IMP in the network grows linearly with N and this too

TABLE II
THE COST OF DISTRIBUTED RESOURCES

Access Method	Collisions	Control Overhead	Idle Capacity
Random access contention	Yes	No	No
Dynamic reservation	No	Yes	No
Fixed allocation	No	No	Yes

places an unacceptable burden on the storage requirements within an IMP. In addition, the transmission and processing costs for updating such large tables is prohibitive. Third, even were the design possible, the cost of the lines connecting this huge number of nodes together grows very quickly unless extreme care is taken in that design. In all three cases just mentioned, one finds that the use of hierarchical structures saves the day. In the design case, one may decompose the network into clusters of nodes, superclusters of clusters, etc., designing each level cluster separately. This significantly reduces the number of nodes involved in each subdesign, thereby reducing the overall design cost significantly. For example, a 100-node net would have a cost on the order of $100^4 = 10^8$ (for $E = 4$), whereas a 2-level hierarchical design with 5 clusters would cost on the order of $5(20)^4 + 5(4) < 10^5$, yielding an improvement of three orders of magnitude! The same approach may be used in routing, where names of distant clusters, rather than names of distant nodes, are used in each routing table, thereby reducing the table length down from N to a number as small as $e \ln N$ giving a significant reduction [24]. For example, a 1000-node net would give a 50-fold reduction in the routing table length when hierarchical routing is used.

In [22] we discuss the overall effect and gain to be had in the use of hierarchically designed wire networks and broadcast networks. For example, we can show that in a bursty dedicated broadcast environment, the use of hierarchical network structures (even with fixed allocation schemes) yields a system cost which is proportional to $(\log M)^2$, where M is the number of users. Comparing this to the case of wire networks where the cost is proportional to the \sqrt{M} , we see the significant advantages that broadcast channels have over wire networks in a bursty environment when hierarchical structures are allowed. We can see this intuitively since we assume that the cost of a broadcast channels is proportional only to capacity, but is independent of distance; if we properly select the transmission range, then the broadcast capacity can be reused spatially (i.e., it can be used independently and simultaneously in more than one area). Further, it can be shown that a 2-level hierarchy using random access in the lower level and dedicated channels in the upper level can be quite efficient in a broadcast environment; this is true since the lower level has gathered together enough traffic so that it is no longer bursty when delivered to the upper level (recall that dedicated channels do well with nonbursty traffic).

V. PRINCIPLES ESTABLISHED

This section is really a continuation of the last since there is a somewhat fuzzy boundary between lessons and principles. Indeed, one might accept the pragmatic definition that a principle is a lesson you had to learn twice.

A. Bigger Is Better

The law of large numbers states that a large collection of demands presents a total demand which is far more predictable

than are the individual demands. We are thus led to the consideration of large shared resources (large in the sense that we increase both the number of users—or the load presented by each user—and the capacity of the resources). Furthermore it is easy to show that the performance improves significantly as we make our systems larger. In particular we can show that a small system whose capacity is C operations per second and whose throughput is J jobs per second (with each job requiring an average of K operations per job) performs A times as slowly (i.e., the response time is A times longer) as a system whose capacity is AC and whose throughput is AJ . The lesson here is very clear, namely that bigger systems perform far better than smaller ones [25]. This is a statement about performance and not one about cost. Indeed if one is talking about communication channel capacity, then one usually also gains through an economy of scale due to the tariff pricing structure as presented by the common carriers. All the more reason, therefore, to concentrate more and more traffic on ever larger channels to gain both cost and efficiency in performance; of course one must be careful not to abuse any "reale" restrictions. Moreover, our lesson about burst communications tells us that in sharing this large channel dynamically, one should provide the full capacity to a single user on demand, rather than to preallocate fractions of the capacity on a permanent basis (omitting consideration of such channel-sharing schemes as spread-spectrum).

The "bigger is better" principle may not apply to the case of stream traffic (defined as real-time traffic which requires a low delay and moderately large throughput requirement—an example being packetized speech). Indeed, an unresolved issue recently raised by Dr. Robert E. Kahn (Editor of this Special Issue) is how effective it would be to handle stream traffic by dividing each trunk into a multiplicity of medium-capacity channels which may then be linked together to form a stream traffic path. We are currently looking at this issue.

B. The Switch

Our second principle has to do with the use of a software switch at the nodes of a network. The principle here is that it pays to place intelligence at the switching nodes of a network since the cost of that intelligence is decreasing far more rapidly than the cost of the communications resource to which it is attached. The idea is to invest some cost in an intelligent switch so as to save yet greater cost in the expensive communications resource. The ability to introduce new programs, new functions, new topologies, new nodes, etc., are all enhanced by the programmable features of a clever communications processor/multiplexer at the software node.

C. Constraints

The principle here is simply, "constraints are necessary and often are evil." Indeed some of the constraints we have seen are sequencing, storage management, capacity allocation, speed matching, and other flow and routing control functions. These "natural" constraints render us vulnerable to dangerous deadlocks and degradations. As mentioned above, if the constraint cannot be met due to some possibly unfortunate accident, then the system will stop all flow. If one is slow in meeting the constraints, then that represents a delay-throughput degradation. As a result of this principle, we see that it behooves us to provide sufficient resources in the network which then allow us to be more relaxed about assigning them. That is, the more precious is a given resource, the tighter we are in

allocating it to a demand, the more likely we are to run into a deadlock or degradation. With an ample resource, we can be more cavalier in assignment and even renege on the assignment if necessary, assuming that a backup facility (in the form of an ample resource elsewhere in the network) is provided.

D. Distributed Control

The principle here is that one must pay a price to nature for organizing a collection of distributed resources into a cooperating group. We have not yet established what that minimum price is, but we have classified the price in the form of collisions, control overhead, and idle capacity.

E. Flow Control

The "principle" here is that flow control is a critical function in packet communications and we are still naive in the invention and analysis of flow control procedures. Hopefully, cleaner code and cleaner concepts will simplify our ability to design and evaluate flow control procedures in the future. There is a "miniprinciple," which seems to be emerging from our preliminary studies [26] which states that if one wants to maximize the power in a network at fixed cost, where power is defined as throughput divided by response time, then under simple statistical assumptions on the flow, one should operate at a point where the throughput delivered is half the maximum possible and the response time is then twice the minimum (no-load) response time.

F. Stale Protocols

In our experiments in the ARPANET, the SATNET, and the packet radio network, we have occasionally attempted to adopt a protocol from one network directly over into a new network. We have found that this is a dangerous procedure and must be carefully analyzed and measured before one adopts such a procedure. Indeed, the use of old protocols in a new environment is dangerous. For example, we found that the use of the ARPANET-like RFNM end-to-end protocol was extremely wasteful of channel capacity and resulted in a capture effect between pairs of users when used in the SATNET. In a 2-user Time Division Multiple Access scheme (in which odd-numbered slots are permanently assigned to user A and even-numbered slots to user B), user A could prevent B from sending any data if he simply started transmitting first in each of his slots since this would require B to devote all of his slots to returning RFNM's to A . Time in the SATNET is divided into fixed length slots (of 30-ms duration). A slot is used for a single packet transmission even if the packet itself is tiny, as is the case of a RFNM. This inefficiency does not exist in the ARPANET since no extra bits are stuffed into ARPANET packets to artificially increase their size. Indeed, gateways have been introduced between the ARPANET and SATNET which renders these nets independent of each other's protocols and formats [2].

G. Inexperienced Designers

It is important that users recognize the difference in function, performance, and operation of a packet network as opposed to a leased line. Certain decisions regarding the parameter settings in any process-to-process communication are often left up to the user of a packet network; for example the buffer allocation he provides in his HOST to accept data from another process communicating through the network with his HOST is a decision often left to the system user. If his buffer allocation is too small, he may degrade the apparent

performance of the network to an unacceptably low degree; this comes about, not because the network is slow, but rather because his allocation was too small. The principle here is that if one leaves design decisions in the hands of the users (or even network designers) then those individuals must be informed as to the effect of their decisions regarding these parameter settings; they cannot be expected to understand the consequence of their actions without being so informed.

VI. CONCLUSIONS

The purpose of this paper has been to boil down a decade of experience with packet communications and from this to extract some lessons and principles we have established. We have succeeded only in part in this endeavor; the field is still moving rapidly and we are learning new things each day. Indeed, in addition to lessons and principles, we have identified a number of open issues which require further study. Aside from the meager principles we stated in the preceding section, we feel it is necessary to make some concluding statements. First we feel that one must view packet communications as a system rather than as a trivial leased line substitute. The use of packet communications offers opportunities to the informed user on the one hand and sets traps for the naive user on the other. It is necessary that the overriding principles which we have established and others which we have yet to establish be well understood by the practitioners in the field. We must continue to learn from our experience, and alas, that experience is often gained through mistakes observed rather than through clever prediction. In all of our design procedures we must constantly be aware of the opportunity to share large resources among large populations of competing demands. We must further be prepared to incorporate new technologies and new applications as they arise; we cannot depend upon "principles" as these principles become invalid in the face of changing technologies and applications.

Lastly, we must point out that the true sharing of processing facilities in the network (i.e., the HOSTS) has not yet been realized in modern day networks. One would dearly love to submit a task to a network, ask that it be accomplished in the most efficient fashion, and expect the network to find the most suitable resources on which to perform that task. Currently, one must specify on which HOST his program should be stored, where his job should be executed, where to store his results, at which location his results should be printed, and specify when all this must happen. The next phase of networking must address this general question of automatic resource sharing among HOSTS in a distributed processing environment. Perhaps in the next special issue on packet communications we will be in a position to identify lessons and principles for true resource sharing of this type.

REFERENCES

- [1] L. Kleinrock, *Queueing Systems, Volume II: Computer Applications*. New York: Wiley Interscience, 1976.
- [2] I. M. Jacobs, R. Binder, and E. V. Hoversten, "General purpose packet satellite networks," this issue, pp. 1448-1467.
- [3] R. E. Kahn, S. A. Gronemeyer, J. Burchfiel, and R. C. Kunzelman, "Advances in packet radio technology," this issue, pp. 1468-1496.
- [4] P. Baran, "On distributed communications," RAND Series Reports, Rand Corporation, Santa Monica, CA, Aug. 1964.
- [5] D. Davies, "The principles of a data communication network for computers and remote peripherals," in *Proc. IFIP Congress '68*, Edinburgh, Scotland, pp. 709-714; Aug. 1968.
- [6] L. G. Roberts, "Multiple computer network development to achieve resource sharing," in *Proc. ACM Symp. Operating Syst.*, Gatlinburg, TN, 1967.
- [7] L. G. Roberts, "The evolution of packet switching," this issue, pp. 1307-1313.
- [8] L. Kleinrock, *Communication Nets; Stochastic Message Flow and Delay*. New York: McGraw-Hill, 1964, out of print. Reprinted by Dover Publications, 1972. (Published in Russian, 1970, Published in Japanese, 1975.)
- [9] O. J. Boxma, "On a tandem queueing model with identical service times at both counters, I," University Utrecht, Dept. of Mathematics, Preprint No. 78, Mar. 1978.
- [10] L. Kleinrock, "Performance of distributed multi-access computer communication systems," in *Proceedings of IFIP Congress '77* Toronto, Canada, pp. 547-552; Aug. 1977.
- [11] L. Pouzin and H. Zimmerman, "A tutorial on protocols," this issue, pp. 1346-1370.
- [12] L. Kleinrock, "ARPANET lessons," in *Proc. Int. Conf. Communications*, Philadelphia, PA, pp. 20-1-20-6; June 1976.
- [13] R. Kahn and W. Crowther, "Flow control in a resource sharing computer network," in *Proc. 2nd IEEE Symp. Problems in Optimization of Data Communication Systems*, Palo Alto, CA, pp. 108-116, Oct. 1971, (also reprinted in *IEEE Trans. Communications*, pp. 539-546; June 1972).
- [14] E. Raubold and J. Haenle, "A method of deadlock-free resource allocation and flow control in packet networks," in *Proc. Third Int. Conf. Computer Communication*, Toronto, Canada, pp. 483-487; Aug. 1976.
- [15] W. Naylor, "A loop-free adaptive routing algorithm for packet switched networks," in *Proc. Fourth Data Communications Symp.*, Quebec City, Canada, pp. 7.9-7.14; Oct. 1975.
- [16] A. Segall, P. M. McKeown, and R. G. Gallager, "A recoverable protocol for loop-free distributed routing," in *Proc. Int. Conf. Communications*, Toronto, Canada, vol. 1, pp. 3.5.1-3.5.5; June 1978.
- [17] H. Opderbeck and L. Kleinrock, "The influence of control procedures on the performance of packet-switched networks," in *Nat. Telecommunications Conf. Record*, San Diego, CA, pp. 810-817; Dec. 1974.
- [18] W. Price, "Simulation of packet-switching networks controlled on isarithmic principles," in *Proc. Third Data Communication Symp.*, St. Petersburg, FL, pp. 44-49; Nov. 1973.
- [19] W. Naylor and L. Kleinrock, "On the effect of periodic routing updates in packet-switched networks," in *Nat. Telecommunications Conf. Record*, Dallas, TX, pp. 16.2-1-16.2-7; Nov. 1976.
- [20] L. Kleinrock and M. Gerla, "On the measured performance of packet satellite access schemes," in *Proc. Fourth Int. Conf. Computer Communication*, Kyoto, Japan, Sept. 1978.
- [21] F. Tobagi, "Performance analysis of packet radio communication systems," in *Nat. Telecommunications Conf. Record*, pp. 12.6-2-12.6-7; Dec. 1977.
- [22] G. Akavia, "Hierarchical organization of distributed packet-switching communication systems," Ph.D. Dissertation, Computer Science Department, Univ. of California, Los Angeles, Mar. 1978.
- [23] R. Metcalfe and D. Boggs, "Ethernet: Distributed packet switching for local computer networks," *Communications of the ACM*, vol. 19, no. 7; pp. 395-404, July 1976.
- [24] L. Kleinrock and F. Kamoun, "Data communications through large packet-switching networks," in *Proc. Int. Teletraffic Congress*, Sydney, Australia, pp. 521-1-521-10; Nov. 1976.
- [25] L. Kleinrock, "Resource allocation in computer systems and computer communication networks," in *Proc. of IFIP Congress '74*, Stockholm, Sweden, pp. 11-18; Aug. 1974.
- [26] —, "On flow control," in *Proc. Int. Conf. Communications*, Toronto, Canada pp. 27.2-1 to 27.2-5, June 1978.
- [27] V. G. T. Cerf and P. Kirstein, "Issues in packet network interconnection," this issue, pp. 1386-1408.